

35.14356

PATENT APPLICATION

IN THE UNITED STATES/PATENT AND TRADEMARK OFFICE

In re Application of:	)	
	:	Examiner: NYA
NORIKI OTANI ET AL.	)	
	:	Group Art Unit: 2776
Application No.: 09/533,255	)	
	:	
Filed: March 23, 2000	)	
	:	
For: APPARATUS AND METHOD	)	
FOR DIVIDING DOCUMENT	:	
INCLUDING TABLE	)	July 27, 2000

Assistant Commissioner for Patents  
Washington, D.C. 20231

CLAIM TO PRIORITY

Sir:

Applicants hereby claims priority under the  
International Convention and all rights to which they are  
entitled under 35 U.S.C. § 119 based upon the following Japanese  
Priority Applications:

11-077583 Filed March 23, 1999  
2000-081870 filed March 23, 2000

Certified copies of the priority documents are  
enclosed.

Applicants' undersigned attorney may be reached in our New York office by telephone at (212) 218-2100. All correspondence should continue to be directed to our address given below.

Respectfully submitted,

  
Attorney for Applicants

Registration No. 25,823

FITZPATRICK, CELLA, HARPER & SCINTO  
30 Rockefeller Plaza  
New York, New York 10112-3801  
Facsimile: (212) 218-2200

NY\_MAIN 99071 v 1

日本国特許庁  
PATENT OFFICE  
JAPANESE GOVERNMENT



別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出願年月日

Date of Application:

1999年 3月23日

出願番号

Application Number:

平成11年特許願第077583号

出願人

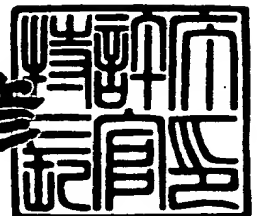
Applicant(s):

キヤノン株式会社

2000年 4月14日

特許庁長官  
Commissioner,  
Patent Office

近藤 隆彦



出証番号 出証特2000-3027058

【書類名】 特許願

【整理番号】 3931016

【提出日】 平成11年 3月23日

【あて先】 特許庁長官 伊佐山 建志 殿

【国際特許分類】 G06F 17/30

【発明の名称】 文書分割装置及び方法、及びそのプログラムを記憶した  
記憶媒体

【請求項の数】 15

【発明者】

【住所又は居所】 東京都大田区下丸子3丁目30番2号キャノン株式会社  
内

【氏名】 大谷 紀子

【発明者】

【住所又は居所】 東京都大田区下丸子3丁目30番2号キャノン株式会社  
内

【氏名】 藤井 憲一

【発明者】

【住所又は居所】 東京都大田区下丸子3丁目30番2号キャノン株式会社  
内

【氏名】 伊藤 史朗

【発明者】

【住所又は居所】 東京都大田区下丸子3丁目30番2号キャノン株式会社  
内

【氏名】 上田 隆也

【発明者】

【住所又は居所】 東京都大田区下丸子3丁目30番2号キャノン株式会社  
内

【氏名】 池田 裕治

【特許出願人】

【識別番号】 000001007  
【住所又は居所】 東京都大田区下丸子 3 丁目 3 0 番 2 号  
【氏名又は名称】 キヤノン株式会社  
【代表者】 御手洗 富士夫  
【電話番号】 03-3758-2111

【代理人】

【識別番号】 100069877  
【住所又は居所】 東京都大田区下丸子 3 丁目 3 0 番 2 号 キヤノン株式会社  
内

【弁理士】

【氏名又は名称】 丸島 儀一  
【電話番号】 03-3758-2111

【手数料の表示】

【予納台帳番号】 011224  
【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1  
【物件名】 図面 1  
【物件名】 要約書 1

【包括委任状番号】 9703271

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 文書分割装置及び方法、及びそのプログラムを記憶した記憶媒体

【特許請求の範囲】

【請求項 1】 処理対象である文書中のテーブルを解析して、各セルの位置関係を示すセル位置データと、各セルの特徴を表現したセルベクトルとを生成するテーブル解析手段と、

該テーブル解析手段により生成されたセル位置データおよびセルベクトルを参照してテーブルタイプを判定するテーブルタイプ判定手段と、

前記テーブルタイプが表を記述したテーブルである場合に、前記テーブルからセグメントを生成する第 1 のセグメント生成手段と、

前記テーブルタイプがレイアウトのためのテーブルである場合に、前記テーブルからセグメントを生成する第 2 のセグメント生成手段とを備えたことを特徴とする文書分割装置。

【請求項 2】 前記第 1 のセグメント生成手段が、

前記セル位置データおよび前記セルベクトルを参照して、前記テーブルにおいて各データが行または列のどちらで表現されているかを判別し、当該テーブルの分割方向を決める分割方向決定手段と、

前記テーブルタイプおよび前記分割方向を参照して、前記テーブルを分割してセグメントを生成する表セグメント生成手段とを備えたことを特徴とする請求項 1 に記載の文書分割装置。

【請求項 3】 前記第 2 のセグメント生成手段が、前記テーブルそのものをセグメントとして生成することを特徴とする請求項 2 に記載の文書分割装置。

【請求項 4】 前記第 2 のセグメント生成手段が、

前記セルベクトルを参照して、前記テーブルにおいて各セルをクラスタリングしてセルクラスタ情報を作成するセルクラスタ作成手段と、

前記セル位置データおよび前記セルクラスタ情報を参照して、前記テーブル中のセルを結合してセグメントを生成するレイアウトセグメント生成手段とを備えたことを特徴とする請求項 1 に記載の文書分割装置。

【請求項 5】 前記第 1 のセグメント生成手段が、前記テーブルそのものをセグメントとして生成することを特徴とする請求項 4 に記載の文書分割装置。

【請求項 6】 前記レイアウトセグメント生成手段が、前記テーブル中のセルの配置パターンを推定し、当該配置パターンに合致するセルを結合してセグメントを生成することを特徴とする請求項 2 に記載の文書分割装置。

【請求項 7】 テーブルを 1 つのセグメントとして文書をセグメントに分割する一般セグメント生成手段を備え、

該一般セグメント生成手段により 1 つのセグメントとして生成されたテーブルを前記テーブル解析手段の処理対象とすることを特徴とする請求項 1 に記載の文書分割装置。

【請求項 8】 処理対象である文書中のテーブルを解析して、各セルの位置関係を示すセル位置データと、各セルの特徴を表現したセルベクトルとを生成するテーブル解析工程と、

該テーブル解析工程により生成されたセル位置データおよびセルベクトルを参照してテーブルタイプを判定するテーブルタイプ判定工程と、

前記テーブルタイプが表を記述したテーブルである場合に、前記テーブルからセグメントを生成する第 1 のセグメント生成工程と、

前記テーブルタイプがレイアウトのためのテーブルである場合に、前記テーブルからセグメントを生成する第 2 のセグメント生成工程とを備えたことを特徴とする文書分割方法。

【請求項 9】 前記第 1 のセグメント生成工程が、

前記セル位置データおよび前記セルベクトルを参照して、前記テーブルにおいて各データが行または列のどちらで表現されているかを判別し、当該テーブルの分割方向を決める分割方向決定工程と、

前記テーブルタイプおよび前記分割方向を参照して、前記テーブルを分割してセグメントを生成する表セグメント生成工程とを備えたことを特徴とする請求項 8 に記載の文書分割方法。

【請求項 10】 前記第 2 のセグメント生成工程が、前記テーブルそのものをセグメントとして生成することを特徴とする請求項 9 に記載の文書分割方法。

【請求項 11】 前記第 2 のセグメント生成工程が、

前記セルベクトルを参照して、前記テーブルにおいて各セルをクラスタリングしてセルクラスタ情報を作成するセルクラスタ作成工程と、

前記セル位置データおよび前記セルクラスタ情報を参照して、前記テーブル中のセルを結合してセグメントを生成するレイアウトセグメント生成工程とを備えたことを特徴とする請求項 8 に記載の文書分割方法。

【請求項 12】 前記第 1 のセグメント生成工程が、前記テーブルそのものをセグメントとして生成することを特徴とする請求項 11 に記載の文書分割方法。

【請求項 13】 前記レイアウトセグメント生成工程が、前記テーブル中のセルの配置パターンを推定し、当該配置パターンに合致するセルを結合してセグメントを生成することを特徴とする請求項 9 に記載の文書分割方法。

【請求項 14】 テーブルを 1 つのセグメントとして文書をセグメントに分割する一般セグメント生成工程を備え、

該一般セグメント生成工程により 1 つのセグメントとして生成されたテーブルを前記テーブル解析工程の処理対象とすることを特徴とする請求項 8 に記載の文書分割方法。

【請求項 15】 処理対象である文書中のテーブルを解析して、各セルの位置関係を示すセル位置データと、各セルの特徴を表現したセルベクトルとを生成するテーブル解析工程と、

該テーブル解析工程により生成されたセル位置データおよびセルベクトルを参照してテーブルタイプを判定するテーブルタイプ判定工程と、

前記テーブルタイプが表を記述したテーブルである場合に、前記テーブルからセグメントを生成する第 1 のセグメント生成工程と、

前記テーブルタイプがレイアウトのためのテーブルである場合に、前記テーブルからセグメントを生成する第 2 のセグメント生成工程とをコンピュータに実行させるための文書分割プログラムを記憶したことを特徴とする記憶媒体。



【発明の詳細な説明】

【0 0 0 1】

【発明の属する技術分野】

本発明は、文書を内容ごとに分割する文書分割装置とその方法、特に、テーブルを含む文書を分割する文書分割装置とその方法に関するものである。

【0 0 0 2】

【従来の技術】

従来、Web上の情報は「ページ」という単位で提供されており、ページの構成や大きさは情報提供者が自由に設定できる。もちろん、情報提供者は各自の情報伝達意図に基づいてページを作成しているのだが、それが必ずしも閲覧者の要求と一致しているとは限らない。

【0 0 0 3】

従って、情報提供者によって関連性が高いと判断された一連の話題が1ページにまとめられていても、閲覧者にとってはそれらの関連性が不要である可能性もあり、複数の話題のうちの1つだけが有用である場合には、他の話題の情報は必要な情報を探索する際の妨げにすらなる。特に、情報提示スペースの小さいモバイル機器では、必要な情報だけを表示するということが重要な機能となる。

【0 0 0 4】

そこで、表示対象である文書をあらかじめ内容ごとに分割しておき、閲覧者が必要としている部分だけを提示することが重要となる。Webページの大半は、Webページ記述言語であるHTML (Hyper Text Markup Language)を用いて書かれている。HTMLは文書構造を記述する言語であるが、論理構造の詳細を記述することは難しく、ブラウザにおけるレイアウトの指定が主な役割となっている。

【0 0 0 5】

しかし、ページのレイアウトには、情報提供者の情報に対する視点が反映されていると考えられる。そこで、情報提供者の意図を反映したセグメントを生成するために、HTMLのタグから読み取ったレイアウトに基づいてページを分割する手法が提案されている。

【0006】

【発明が解決しようとする課題】

上記提案の手法では、<TABLE>タグと</TABLE>タグで囲まれたテーブルは、意味的なまとまりであると判断されて、1つのセグメントとして形成されている。しかしながら、テーブルは、比較的大きな領域を占めて複数の情報を含んでいる場合が多いため、さらに細かいセグメントに分割することが望ましい。

【0007】

その際、テーブルは、単純な表を記述している場合と、テキストやイメージのレイアウトを指定している場合とがあるが、両者においてタグに含まれた意図はまったく異なるので、それぞれ違うアプローチでセグメントを生成すべきである。

【0008】

単純な表を記述している場合は、含まれているデータごとにセグメントを生成することで、ユーザのより細かい要求に備えることができると考えられる。ところが、一口に表を記述していると言っても、1組のデータが行で表現されていたり列で表現されていたり、項目名を記述した行(または列)があったりなかったりと、様々な表の形式が存在する。従って、表をデータごとのセグメントに分割するためには、表の形式を判断する必要がある。

【0009】

一方、テキストやイメージをレイアウトするためにテーブルタグを使っている場合は、各セルに記述された内容とセル同士の位置関係からセル間の関係を推定し、内容のまとまりを判断してセグメントを生成することが望まれる。

【0010】

本発明は、上記の課題に鑑みてなされたものであり、処理対象となっているテーブルを解析して、表を記述したテーブルであるか、レイアウト目的のテーブルであるかを判別し、それぞれに応じた処理によってセグメントを生成することで、文書中のテーブルを内容ごとに分割する文書分割装置を提供することを目的とする。

## 【0011】

## 【課題を解決するための手段】

上述した目的を達成するために、本発明によれば、文書分割装置に、処理対象である文書中のテーブルを解析して、各セルの位置関係を示すセル位置データと、各セルの特徴を表現したセルベクトルとを生成するテーブル解析手段と、該テーブル解析手段により生成されたセル位置データおよびセルベクトルを参照してテーブルタイプを判定するテーブルタイプ判定手段と、前記テーブルタイプが表を記述したテーブルである場合に、前記テーブルからセグメントを生成する第1のセグメント生成手段と、前記テーブルタイプがレイアウトのためのテーブルである場合に、前記テーブルからセグメントを生成する第2のセグメント生成手段とを備える。

## 【0012】

また、本発明の他の態様によれば、文書分割方法に、処理対象である文書中のテーブルを解析して、各セルの位置関係を示すセル位置データと、各セルの特徴を表現したセルベクトルとを生成するテーブル解析工程と、該テーブル解析工程により生成されたセル位置データおよびセルベクトルを参照してテーブルタイプを判定するテーブルタイプ判定工程と、前記テーブルタイプが表を記述したテーブルである場合に、前記テーブルからセグメントを生成する第1のセグメント生成工程と、前記テーブルタイプがレイアウトのためのテーブルである場合に、前記テーブルからセグメントを生成する第2のセグメント生成工程とを備える。

## 【0013】

更に、本発明の他の態様によれば、記憶媒体に、処理対象である文書中のテーブルを解析して、各セルの位置関係を示すセル位置データと、各セルの特徴を表現したセルベクトルとを生成するテーブル解析工程と、該テーブル解析工程により生成されたセル位置データおよびセルベクトルを参照してテーブルタイプを判定するテーブルタイプ判定工程と、前記テーブルタイプが表を記述したテーブルである場合に、前記テーブルからセグメントを生成する第1のセグメント生成工程と、前記テーブルタイプがレイアウトのためのテーブルである場合に、前記テーブルからセグメントを生成する第2のセグメント生成工程とをコンピュータに

実行させるための文書分割プログラムを記憶する。

【0014】

【発明の実施の形態】

以下、図面を用いて本発明の1実施形態を詳細に説明する。

【0015】

図1は、本実施形態の文書分割装置の機能構成を示すブロック図である。同図において、101は、処理対象であるHTML文書中のテーブル(<table>と</table>で囲まれた部分)を保持するHTMLテーブル保持部である。

【0016】

102は、HTMLテーブル保持部101に保持されているテーブルを解析して、各セルの位置関係を示すセル位置データと、各セルの特徴を表現したセルベクトルとを生成するテーブル解析部である。

【0017】

セルベクトルは、セルの高さや幅、内容の表示位置、背景色、セル内のテキストの長さや文字種、セル内のイメージの大きさや形などから決定する。セルベクトルの次元は(セル内のイメージの個数×4+17)次元であり、各成分は0以上1以下の実数である。

【0018】

セル内でi番目に出現するイメージを $image_i$ とするとき、セルベクトル $v$ の第 $k$ 成分 $v(k)$ は次のように定義される。

【0019】

$v(0)$  : タグの種類が<TH>(項目名を表現するセル)のとき1.0、<TD>(データを表現するセル)のとき0.0。

【0020】

$v(1)$  : rowspan (行幅) が4未満のとき $rowspan \times 0.25$ 、4以上のとき1.0。

【0021】

$v(2)$  : colspan (列幅) が4未満のとき $colspan \times 0.25$ 、4以上のとき1.0。

【0022】

$v(3)$  : nowrap (改行なし) が指定されているとき1.0、指定されていないとき

0.0。

【0 0 2 3】

v(4) : align (横位置) の指定がないとき0.0、left (左詰め) のとき0.2、center (中央) のとき0.4、right (右詰め) のとき0.6、justify (均等) のとき0.8、それ以外るとき1.0。

【0 0 2 4】

v(5) : valign (縦位置) の指定がないとき0.0、top (上詰め) のとき0.2、middle (中央) のとき0.4、bottom (下詰め) のとき0.6、baseline (ベースライン) のとき0.8、それ以外るとき1.0。

【0 0 2 5】

v(6) : bgcolor (背景色) の指定がないとき0.0、16進コードで指定されていないとき0.0、16進コードで指定されているときbgcolor/0xFFFFFF。

【0 0 2 6】

v(7) : 9列目以前るとき(列番号)×0.1、10列目以降るとき1.0。

【0 0 2 7】

v(8) : 99行目以前るとき(行番号)×0.01、100行目以降るとき1.0。

【0 0 2 8】

v(9) : 改行(<BR>)数が5つ未満るとき(改行数)×0.2、5つ以上るとき1.0。

【0 0 2 9】

v(10) : テキストの文字数が100文字未満るとき(文字数)×0.01、100文字以上るとき1.0。

【0 0 3 0】

v(11) : (テキスト中の数字の数)/(テキストの全文字数)。

【0 0 3 1】

v(12) : (テキスト中のアルファベットの数)/(テキストの全文字数)。

【0 0 3 2】

v(13) : (テキスト中の漢字の数)/(テキストの全文字数)。

【0 0 3 3】

v(14) : (テキスト中のカタカナの数)/(テキストの全文字数)。

【 0 0 3 4 】

v(15) : (テキスト中のひらがなの数)/(テキストの全文字数)。

【 0 0 3 5 】

v(16) : 句点(“。” または “.”)があるとき1.0、ないとき0.0。

【 0 0 3 6 】

v(13+i×4) : image<sub>i</sub>の面積が150000未満のとき(面積)/150000、150000以上のとき1.0。

【 0 0 3 7 】

v(14+i×4) : image<sub>i</sub>の高さが300未満のとき(高さ)/300、300以上のとき1.0。

【 0 0 3 8 】

v(15+i×4) : image<sub>i</sub>の幅が500未満のとき(幅)/500、500以上のとき1.0。

【 0 0 3 9 】

v(16+i×4) : このテーブルを含んでいるページのURLを示す文字列のうち、image<sub>i</sub>のURLと共通の部分文字列の割合。例えば、

http://hoge hoge.aaa.bbbbbb.co.jp:8080/hoge1/hoge2/hoge.html

のページ(URLの長さは58)に“../image/hoge.gif”というイメージがあった場合、イメージをフルパスのURLに書き換えると、

http://hoge hoge.aaa.bbbbbb.co.jp:8080/hoge1/image/hoge.gif

となるので、共通の部分文字列は

http://hoge hoge.aaa.bbbbbb.co.jp:8080/hoge1/

となる。この長さは43なので、この成分の値は43÷58=0.741となる。

【 0 0 4 0 】

103は、テーブル解析部102により生成されたセル位置データを保持するセル位置データ保持部である。104は、テーブル解析部102により生成されたセルベクトルを保持するセルベクトル保持部である。

【 0 0 4 1 】

105は、セル位置データ保持部103に保持されたセル位置データ、およびセルベクトル保持部104に保持されたセルベクトルを参照してテーブルタイプを判定し

、テーブルタイプによってカット方向決定部107、またはセルクラス作成部111に処理開始を指示するテーブルタイプ判定部である。テーブルタイプには、以下のtable I～table VIIの7種類がある。

【0 0 4 2】

table I： すべてのセルの高さと幅が1であり、1行n列目及びn行1列目のセルがすべて<TH>または同じ背景色。

【0 0 4 3】

table II： すべてのセルの高さと幅が1であり、1行n列目及びn行1列目(1行1列目を除く)のセルがすべて<TH>または同じ背景色。

【0 0 4 4】

table III： すべてのセルの高さと幅が1であり、1行n列目のセルがすべて<TH>または同じ背景色。

【0 0 4 5】

table IV： すべてのセルの高さと幅が1であり、1行n列目(1行1列目を除く)のセルがすべて<TH>または同じ背景色。

【0 0 4 6】

table V： すべてのセルの高さと幅が1であり、n行1列目のセルがすべて<TH>または同じ背景色。

【0 0 4 7】

table VI： すべてのセルの高さと幅が1であり、n行1列目(1行1列目を除く)のセルがすべて<TH>または同じ背景色。

【0 0 4 8】

table VII： table I～table VI以外のテーブル。

【0 0 4 9】

以上において、table I～table VIが表を記述するためのテーブルであり、table VIIがレイアウト目的のテーブルである。テーブルタイプがtable I～table V Iの場合にはカット方向決定部107に処理開始を指示し、テーブルタイプがtable VIIの場合にはセルクラス作成部111に処理開始を指示する。

## 【 0 0 5 0 】

106は、テーブルタイプ判定部105により決定されたテーブルタイプを保持するテーブルタイプ保持部である。

## 【 0 0 5 1 】

107は、テーブルタイプ判定部105により処理開始を指示された場合に、セル位置データ保持部103に保持されたセル位置データ、およびセルベクトル保持部104に保持されたセルベクトルを参照して、表を記述したテーブルにおいて各データは行または列のどちらで表現されているかを判別し、テーブルの分割方向を決めるカット方向決定部である。

## 【 0 0 5 2 】

N行M列のテーブルTを行で分割したときのスコア $S_h(T)$ と列で分割したときのスコア $S_v(T)$ を以下のように定義する。以下で、 $\cos(v_{i,j}, v_{k,l})$ はi行j列目のセルのテーブルセルベクトル $v_{i,j}$ とk行l列目のセルのテーブルセルベクトル $v_{k,l}$ との余弦値を表す。

## 【 0 0 5 3 】

ただし、これはi行j列目のセルとk行l列目のセルのデータとが共に存在するときのみ算出される値で、両方もしくはどちらか一方のセルのデータが存在しない場合には、値は0となる。



【 0 0 5 4 】

【外 1】

$$\text{exist}(i,j) = \begin{cases} 1 & (i \text{ 行 } j \text{ 目列のセルにデータが存在する}) \\ 0 & (i \text{ 行 } j \text{ 目列のセルにデータが存在しない}) \end{cases}$$

$$\text{count}_h = \sum_{i=1}^N \sum_{j=1}^M \sum_{k=j+1}^M \text{exist}(i,j) \times \text{exist}(i,k)$$

$$\text{count}_v = \sum_{j=1}^M \sum_{i=1}^N \sum_{l=j+1}^M \text{exist}(i,j) \times \text{exist}(i,l)$$

$$S_h(T) = \frac{1}{\text{count}_h} \sum_{i=1}^N \sum_{j=1}^M \sum_{k=j+1}^M \cos(v_{i,j}, v_{i,k})$$

$$S_v(T) = \frac{1}{\text{count}_v} \sum_{j=1}^M \sum_{i=1}^N \sum_{l=j+1}^M \cos(v_{i,j}, v_{i,l})$$

【 0 0 5 5 】

テーブルセルベクトルの次元は、i行j列目のセルとk行l列目のセルに含まれるイメージの数により決定されるので、両ベクトルの次元が同じになるように、低次元のテーブルセルベクトルに値0の成分を追加して余弦値を計算する。

【 0 0 5 6 】

$S_h(T)$ は同じ行にある2つのセルのテーブルセルベクトルの平均余弦値であり、 $S_v(T)$ は同じ列にある2つのセルのテーブルセルベクトルの平均余弦値である。2つのテーブルセルベクトルの余弦値はセルの類似度と見なせるので、 $S_h(T)$ はテーブルを行ごとに分割した時の同セグメント内におけるセル間の平均類似度、 $S_v(T)$ はテーブルを列ごとに分割した時の同セグメント内におけるセル間の平均類似度といえる。

【 0 0 5 7 】

各セグメントに各種のデータを盛り込むには、同セグメント内セル間類似度が低い方が良いので、 $S_h(T) \leq S_v(T)$ のときはテーブルTを行ごとに分割し、 $S_h(T) > S_v(T)$ のときテーブルTを列ごとに分割するべきだと判断する。

【 0 0 5 8 】

108は、カット方向決定部107により決定されたカット方向を保持するカット方

向保持部である。

【0059】

109は、テーブルタイプ保持部106に保持されたテーブルタイプ、およびカット方向保持部108に保持されたカット方向を参照して、表を記述したテーブルからセグメントを生成する表セグメント生成部である。カット方向が行方向の場合、table Vのテーブルはそのまま行をセグメントとし、table V以外のテーブルは1行目を組み合わせてセグメントを作る。カット方向が列方向の場合、table IIIのテーブルはそのまま列をセグメントとし、table III以外のテーブルは1列目を組み合わせてセグメントを作る。

【0060】

110は、表セグメント生成部109により生成された表セグメントを保持する表セグメント保持部である。

【0061】

111は、テーブルタイプ判定部105により処理開始を指示された場合に、セルベクトル保持部104に保持されたセルベクトルを参照して、レイアウト目的のテーブルにおいて各セルをクラスタリングするセルクラスタ作成部である。ここでは最大距離アルゴリズムを用いてセルの分類を決定する。最大距離アルゴリズムのクラスタリング手順を以下に示す。

【0062】

Step.1: N個のサンプルパターン集合 $X = \{x_1, x_2, \dots, x_N\}$  から、任意にひとつ(ここでは $x_1$ として説明する)を選び、クラスタ中心 $z_1 \in Z$ とする。

【0063】

Step.2: Zに含まれないすべての $x_i \in X$ について、すでに選ばれたクラスタ中心 $z_j \in Z$ のうち、一番近いものまでの距離 $dx_i$ を計算する。 $\text{Max}\{dx_i\}$ を与える $x_i$ を $x_C$ とする。

【0064】

Step.3: すべての $z_k \in Z$ について、 $z_k$ 以外のクラスタ中心のうち、一番遠いものまでの距離 $dz_k$ を計算する。

【 0 0 6 5 】

Step.4:  $dx_C \geq \max \{dz_k\} \times t (t=0.5 \sim 1)$  が成立するとき、 $x_C$  を新たなクラスタ中心とし、Step.2に戻って次のクラスタ中心を選ぶ。 $dx_C < \max \{dz_k\} \times t (t=0.5 \sim 1)$  ならばStep.5へ。

【 0 0 6 6 】

Step.5: すべての  $x_i \in X$  を、最も近い  $z_j \in Z$  のクラスタに分類する。

【 0 0 6 7 】

最大距離アルゴリズムによるクラスタリング結果の例を図4に示す。

【 0 0 6 8 】

112は、セルクラスタ作成部111により作成されたセルのクラスタ情報を保持するセルクラスタ情報保持部である。

【 0 0 6 9 】

113は、セル位置データ保持部103に保持されたセル位置データ、およびセルクラスタ情報保持部112に保持されたセルクラスタ情報を参照して、レイアウト目的のテーブルからセグメントを生成するレイアウトセグメント生成部である。

【 0 0 7 0 】

テーブルの形式を利用して情報を配置するメリットとしては、ある配置パターンの縦横方向の繰り返しが容易に表現できる点が挙げられる。そこで、セルクラスタ情報をもとに配置パターンを推定して、パターンに適合するセルを合わせてセグメントとする。ある配置パターンが繰り返し現れるときには、そのパターンに適合するセル同士が意味的にまとまっていると判断できるからである。処理の詳細を以下に示す。

【 0 0 7 1 】

まず、基本セル種を決定し、基本セル種に属するセルを基本セルとする。基本セル種は、同種のセルが複数あるセルの種類のうち、最もセル数の少ないセル種とする。該当するセル種が複数ある場合には、より左、上にあるセルの種類を選ぶ。

【 0 0 7 2 】

次に、ある基本セルに隣接するセルと分類が同じセルが他の基本セルにも同じ

ように隣接するかを確認する。隣接していれば、それぞれを結合し、新たな基本セルとする。これを結合できなくなるまで繰り返す。

【 0 0 7 3 】

以上の処理を終えると、基本セルおよび残りのセルをそれぞれセグメントとする。

【 0 0 7 4 】

114は、レイアウトセグメント生成部113により生成されたレイアウトセグメントを保持するレイアウトセグメント保持部である。表セグメント保持部110に保持された表セグメント、およびレイアウトセグメント保持部114に保持されたレイアウトセグメントが結果として得られるセグメントである。

【 0 0 7 5 】

図 2 は、本発明の実施形態に係る文書分割装置のハードウェア構成を示す図である。

【 0 0 7 6 】

同図において、201は後述する制御手順を実現するプログラムを保持するROMである。202はRAMで、セル位置データ保持部103、セルベクトル保持部104、テーブルタイプ保持部106、カット方向保持部108、セルクラス情報保持部112と上記プログラムの動作に必要な記憶領域とを提供する。

【 0 0 7 7 】

203はROM201に保持されているプログラムに従って処理を行なう中央処理装置である。204はディスク装置であり、HTMLテーブル保持部101、表セグメント保持部110、レイアウトセグメント保持部114を実現する。205はバスであり、上記の各構成を接続し、各構成間におけるデータの授受を可能とする。

【 0 0 7 8 】

次に、本実施形態の処理動作を説明する。図 3 は本実施形態の文書分割装置の動作手順を示すフローチャートである。

【 0 0 7 9 】

ステップS301では、HTMLテーブル保持部101に保持されているテーブルを解析して、各セルの位置関係を示すセル位置データと、各セルの特徴を表現したセル

ベクトルを生成する。そしてステップS302に移る。

【0080】

ステップS302では、セル位置データ保持部103に保持されたセル位置データ、およびセルベクトル保持部104に保持されたセルベクトルを参照して、テーブルタイプを決定する。そしてステップS303に移る。

【0081】

ステップS303では、テーブルタイプ保持部106に保持されたテーブルタイプを参照して、処理対象のテーブルが表を記述したテーブルか否かを判定して、表を記述したテーブルの場合はステップS304に移る。表を記述したテーブルでない場合はステップS306に移る。

【0082】

ステップS304では、セル位置データ保持部103に保持されたセル位置データ、およびセルベクトル保持部104に保持されたセルベクトルを参照して、表を記述したテーブルにおいて各データは行または列のどちらで表現されているかを判別し、テーブルの分割方向を決める。そしてステップS305に移る。

【0083】

ステップS305では、テーブルタイプ保持部106に保持されたテーブルタイプ、およびカット方向保持部108に保持されたカット方向を参照して、表を記述したテーブルからセグメントを生成する。そして動作を終了する。

【0084】

ステップS306では、セルベクトル保持部104に保持されたセルベクトルを参照して、レイアウト目的のテーブルにおいて各セルをクラスタリングする。そしてステップS307に移る。

【0085】

ステップS307では、セル位置データ保持部103に保持されたセル位置データ、およびセルクラスタ情報保持部112に保持されたセルクラスタ情報を参照して、レイアウト目的のテーブルからセグメントを生成する。そして動作を終了する。

【0086】

以上に述べたように、処理対象となっているテーブルを解析して、表を記述し

たテーブルであるか、レイアウト目的のテーブルであるかを判別し、それぞれに応じた処理によってセグメントを生成することで、HTML文書中のテーブルを内容ごとに分割する文書分割装置を実現することができる。

【0087】

〔他の実施形態〕

上記実施形態では、セルのクラスタリングに最大距離アルゴリズムを利用するよう説明しているが、これに限定されるものではなく、他のアルゴリズムを用いてクラスタリングを行なってもよい。

【0088】

上記実施形態で示したセルベクトルの各成分の定義は一例であり、他の定義によってセルの特徴をベクトル表現してもよい。

【0089】

上記実施形態で示したカット方向を決定するスコアの定義は一例であり、他の定義によってカット方向を決定してもよい。

【0090】

上記実施形態では、テーブルタイプを決定するための項目名の行(または列)の判定に、セルの高さと幅、タグの種類(TH or TD)、背景色を用いているが、これに限定されるものではなく、他の属性を用いて判定してもよい。

【0091】

上記実施形態では、HTMLのテーブルを分割するだけの装置として説明しているが、これに限定されるものではない。例えば、HTML文書全体を分割する装置であってもよい。図5は、この場合の基本的な機能構成を示すブロック図である。

【0092】

図5において、501は、処理対象であるHTML文書を保持するHTML文書保持部である。502は、HTML文書保持部501に保持されているHTML文書をセグメントに分割する一般セグメント生成部である。503は、一般セグメント生成部502により生成されたテーブル以外のセグメントを保持する一般セグメント保持部である。504は、一般セグメント生成部502により生成されたテーブルのセグメントを保持するHTMLテーブル保持部である。

【 0 0 9 3 】

以下、505～517は、図1の102～114と同様である。

【 0 0 9 4 】

図5では、一般セグメント保持部503に保持された一般セグメント、表セグメント保持部513に保持された表セグメント、およびレイアウトセグメント保持部517に保持されたレイアウトセグメントが結果として得られるセグメントである。

【 0 0 9 5 】

上記実施形態では、表を記述しているテーブルとレイアウト目的のテーブルの両方をセグメントに分割しているが、これに限定されるものではない。例えば、表を記述しているテーブルのみを分割してもよい。図6はこの場合の基本的な機能構成を示すブロック図である。

【 0 0 9 6 】

図6において、601～610は、図1の101～110と同様である。

【 0 0 9 7 】

611は、テーブルタイプ判定部605により処理開始を指示された場合に、HTMLテーブル保持部601に保持されたHTMLテーブルをテーブルセグメントとするテーブルセグメント生成部である。

【 0 0 9 8 】

612は、テーブルセグメント生成部611により生成されたテーブルセグメントを保持するテーブルセグメント保持部である。

【 0 0 9 9 】

図6では、表セグメント保持部610に保持された表セグメント、およびテーブルセグメント保持部612に保持されたテーブルセグメントが結果として得られるセグメントである。

【 0 1 0 0 】

また、上記実施形態では、表を記述しているテーブルとレイアウト目的のテーブルの両方をセグメントに分割しているが、レイアウト目的のテーブルのみを分割してもよい。図7はこの場合の基本的な機能構成を示すブロック図である。

【 0 1 0 1 】

図 7 において、701～705 及び 708～711 は、図 1 の 101～115 及び 111～114 と同様である。

【 0 1 0 2 】

706 は、テーブルタイプ判定部 705 により処理開始を指示された場合に、HTML テーブル保持部 701 に保持された HTML テーブルをテーブルセグメントとするテーブルセグメント生成部である。707 は、テーブルセグメント生成部 706 により生成されたテーブルセグメントを保持するテーブルセグメント保持部である。

【 0 1 0 3 】

図 7 では、テーブルセグメント保持部 707 に保持されたテーブルセグメント、およびレイアウトセグメント保持部 711 に保持されたレイアウトセグメントが、結果として得られるセグメントである。

【 0 1 0 4 】

上記実施形態では、HTML 文書を分割する装置として説明しているが、これに限定されるものではなく、検索装置と組み合わせて、生成されたセグメント単位で検索を行なうことができるセグメント検索装置として実現してもよい。

【 0 1 0 5 】

上記実施形態においては、セル位置データ保持部 103、セルベクトル保持部 104、テーブルタイプ保持部 106、カット方向保持部 108、セルクラスタ情報保持部 112 を RAM で、HTML テーブル保持部 101、表セグメント保持部 110、レイアウトセグメント保持部 114 をディスク装置で実現する場合について説明したが、これに限定されるものではなく、任意の記憶媒体を用いて実現してもよい。

【 0 1 0 6 】

上記実施形態では、HTML のテーブルを分割する場合について説明したが、テーブルの内容が区別できれば、他の形式であってもよい。

【 0 1 0 7 】

上記実施形態においては、各部を同一の計算機上で構成する場合について説明したが、これに限定されるものではなく、ネットワーク上に分散した計算機や処理装置などに分かれて各部を構成してもよい。



## 【0108】

上記実施形態においては、プログラムをROMに保持する場合について説明したが、これに限定されるものではなく、任意の記憶媒体を用いて実現してもよい。また、同様の動作をする回路で実現してもよい。

## 【0109】

なお、本発明は、複数の機器から構成されるシステムに適用しても、1つの機器からなる装置に適用してもよい。前述した実施形態の機能を実現するソフトウェアのプログラムコードを記録した記録媒体を、システム或いは装置に供給し、そのシステム或いは装置のコンピュータ（またはCPUやMPU）が記録媒体に格納されたプログラムコードを読み出し実行することによっても、達成されることは言うまでもない。

## 【0110】

この場合、記録媒体から読み出されたプログラムコード自体が前述した実施形態の機能を実現することになり、そのプログラムコードを記録した記録媒体は本発明を構成することになる。

## 【0111】

プログラムコードを供給するための記録媒体としては、例えば、フロッピーディスク、ハードディスク、光ディスク、光磁気ディスク、CD-ROM、CD-R、磁気テープ、不揮発性のメモ리카ード、ROMなどを用いることができる。

## 【0112】

また、コンピュータが読み出したプログラムコードを実行することにより、前述した実施形態の機能が実現されるだけでなく、そのプログラムコードの指示に基づき、コンピュータ上で稼働しているOSなどが実際の処理の一部または全部を行ない、その処理によって前述した実施形態の機能が実現される場合も含まれることは言うまでもない。

## 【0113】

更に、記録媒体から読み出されたプログラムコードが、コンピュータに挿入された機能拡張ボードやコンピュータに接続された機能拡張ユニットに備わるメモリに書き込まれた後、そのプログラムコードの指示に基づき、その機能拡張ボ

ードや機能拡張ユニットに備わるCPUなどが実際の処理の一部または全部を行ない、その処理によって前述した実施形態の機能が実現される場合も含まれることは言うまでもない。

#### 【0114】

#### 【発明の効果】

以上説明したように、本発明によれば、文書中のテーブルを解析して、各セルの位置関係を示すセル位置データと、各セルの特徴を表現したセルベクトルとを生成し、このセル位置データおよびセルベクトルを参照して、処理対象のテーブルが表を記述したテーブルか否かを判定し、判定結果に応じた手法でセグメントを生成することで、文書中のテーブルを内容ごとに分割する文書分割を実現できるという効果が得られる。

#### 【図面の簡単な説明】

#### 【図1】

本発明に係る一実施形態の文書分割装置の基本構成を示すブロック図である。

#### 【図2】

実施形態に係る文書分割装置のハードウェア構成を示すブロック図である。

#### 【図3】

実施形態に係る文書分割装置の動作手順を示すフローチャートである。

#### 【図4】

最大距離アルゴリズムを説明する図である。

#### 【図5】

他の実施形態の基本構成を示すブロック図である。

#### 【図6】

他の実施形態の基本構成を示すブロック図である。

#### 【図7】

他の実施形態の基本構成を示すブロック図である。

#### 【符号の説明】

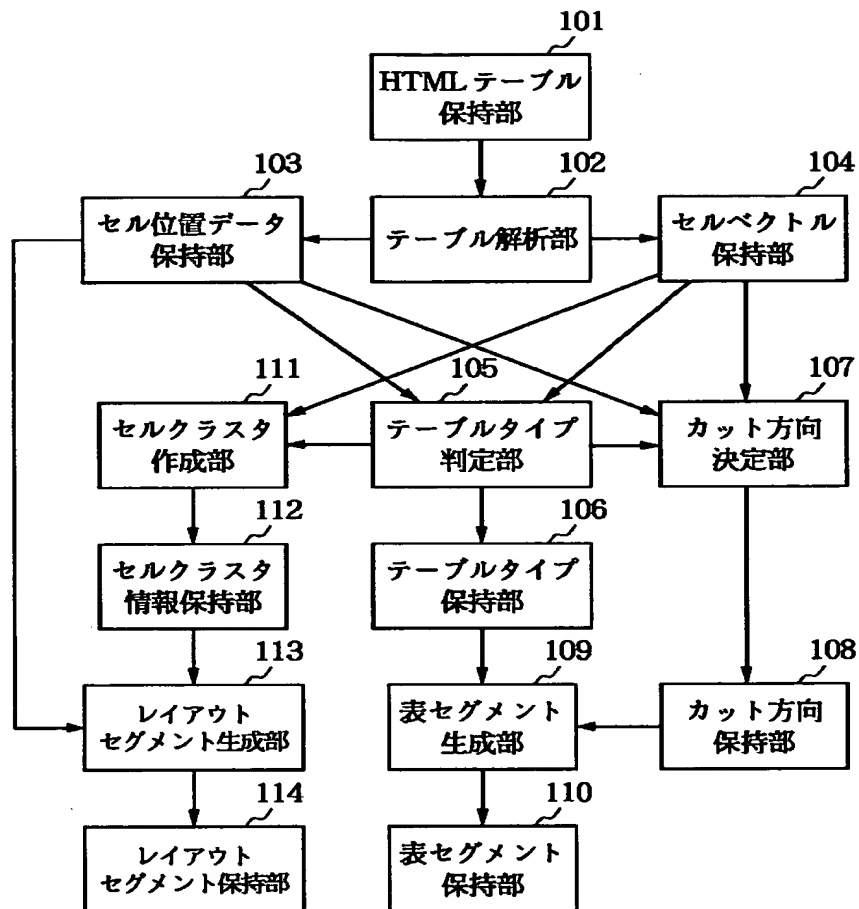
101 HTMLテーブル保持部

102 テーブル解析部

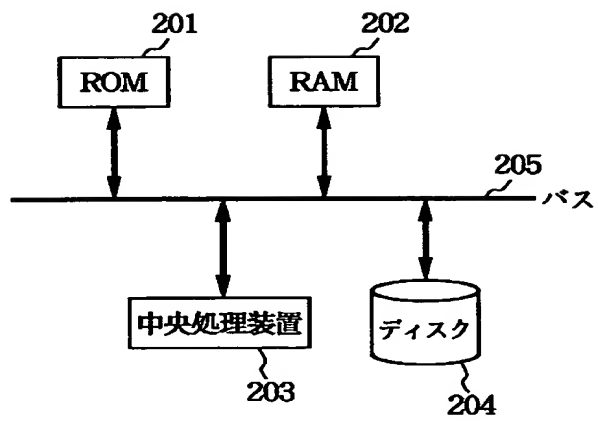
- 1 0 3 セル位置データ保持部
- 1 0 4 セルベクトル保持部
- 1 0 5 テーブルタイプ判定部
- 1 0 6 テーブルタイプ保持部
- 1 0 7 カット方向決定部
- 1 0 8 カット方向保持部
- 1 0 9 表セグメント生成部
- 1 1 0 表セグメント保持部
- 1 1 1 セルクラスタ作成部
- 1 1 2 セルクラスタ情報保持部
- 1 1 3 レイアウトセグメント生成部
- 1 1 4 レイアウトセグメント保持部
- 2 0 1 R O M
- 2 0 2 R A M
- 2 0 3 中央処理装置
- 2 0 4 ディスク
- 2 0 5 バス

【書類名】 図面

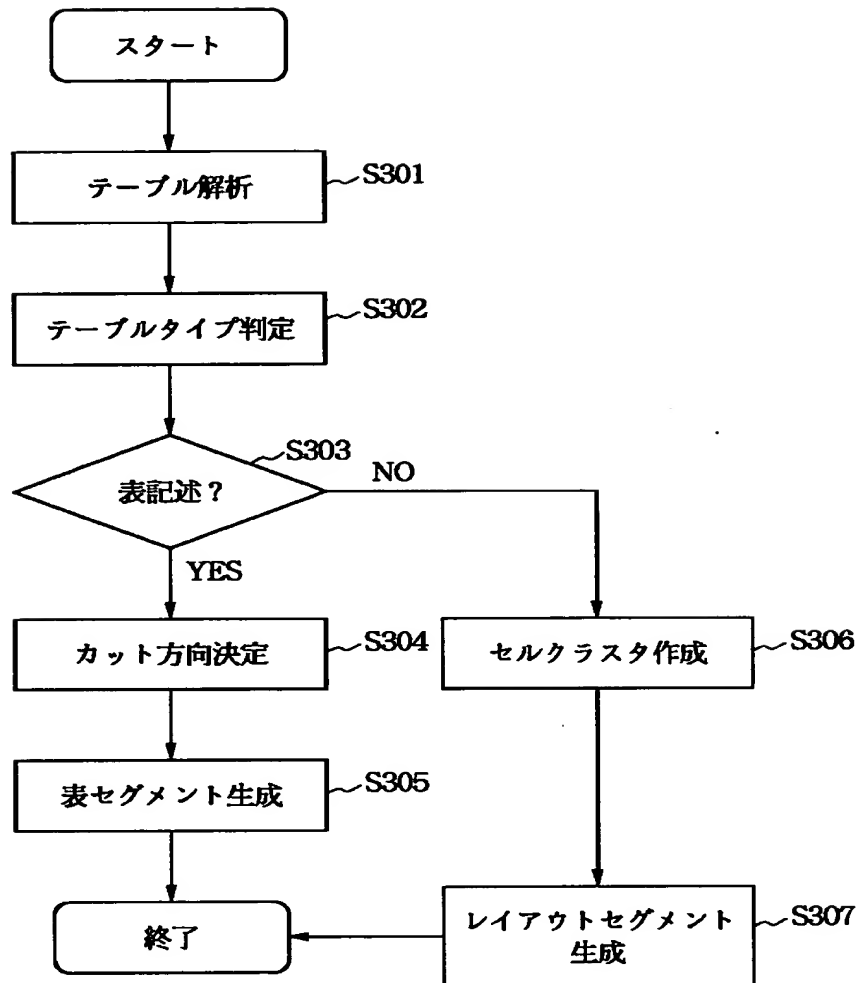
【図 1】



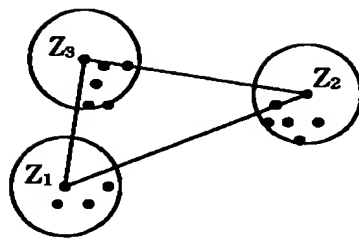
【図 2】



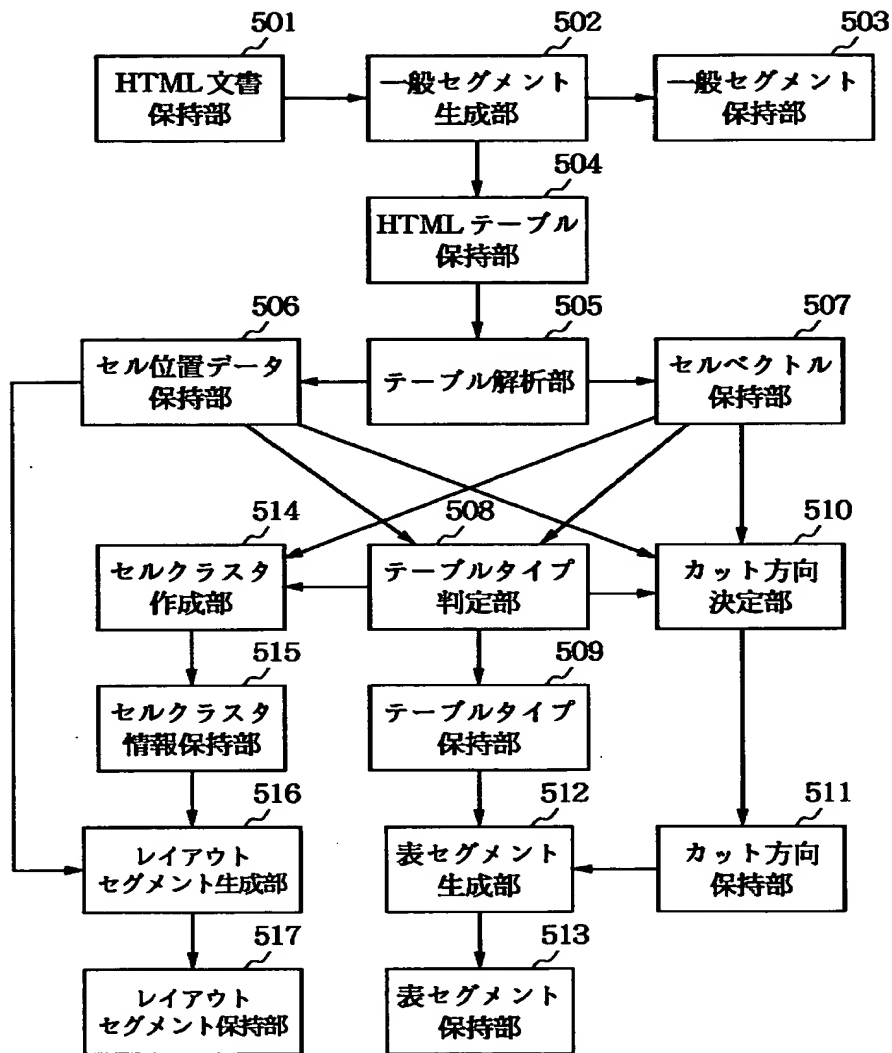
【図 3】



【図 4】

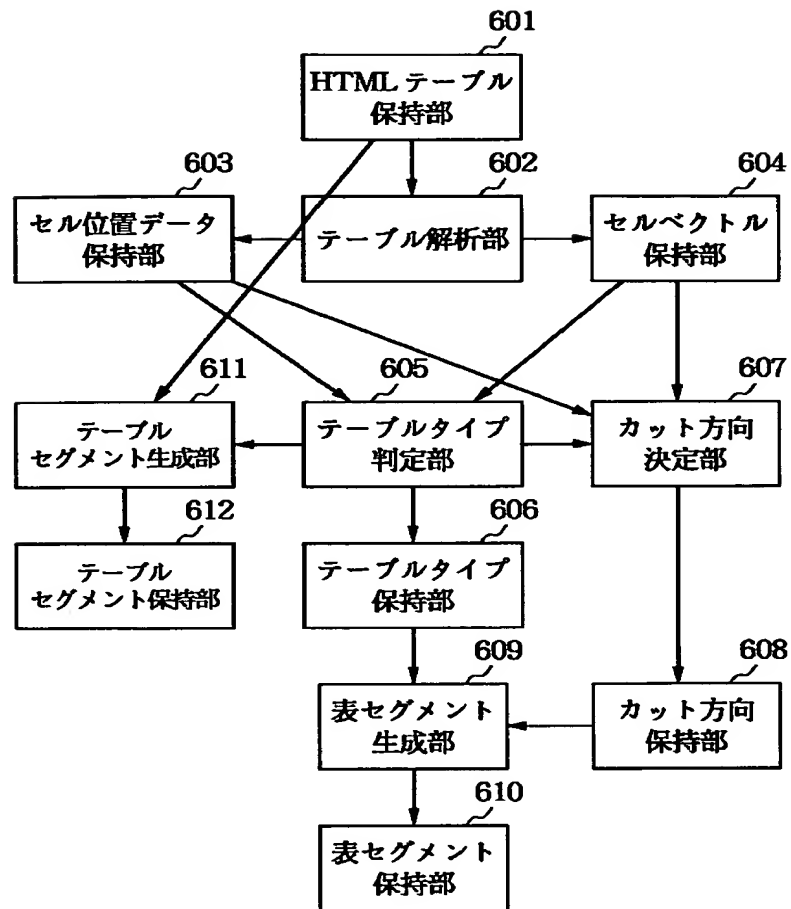


【図 5】

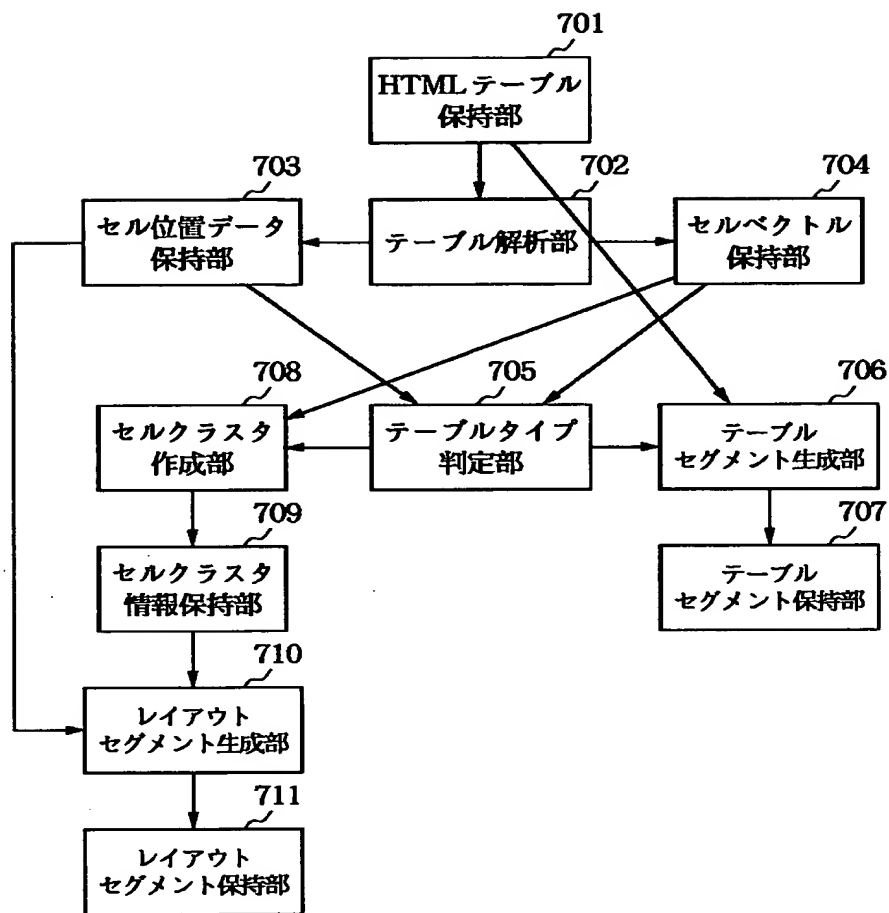




【図 6】



【図 7】



【書類名】 要約書

【要約】

【課題】 HTML文書中のテーブルを内容ごとに分割する。

【解決手段】 HTML文書中のテーブルを解析して、各セルの位置関係を示すセル位置データと各セルの特徴を表現したセルベクトルとを生成し(S301)、このセル位置データおよびセルベクトルを参照してテーブルタイプを判定し(S302)、表を記述したテーブルの場合は、セル位置データおよびセルベクトルを参照して、各データは行または列のどちらで表現されているかを判別し、テーブルの分割方向を決め(S304)、テーブルタイプおよび分割方向を参照してセグメントを生成し(S305)、表を記述したテーブルでないレイアウト目的のテーブルの場合は、セルベクトルを参照して各セルをクラスタリングし(S306)、セル位置データおよびセルクラスタ情報を参照してセグメントを生成する(S307)。

【選択図】 図 3

出 願 人 履 歴 情 報

識別番号 [000001007]

1. 変更年月日	1990年 8月30日
[変更理由]	新規登録
住 所	東京都大田区下丸子3丁目30番2号
氏 名	キヤノン株式会社